

An Iterative Algorithm for Computing the Optimal Exponent of Correct Decoding Probability for Rates below the Rate Distortion Function

Yutaka Jitsumatsu
Dept. of Informatics,
Kyushu University, Japan
Email: jitumatu@inf.kyushu-u.ac.jp

Yasutada Oohama
Dept. of Communication Engineering and Informatics,
University of Electro-Communication, Japan
Email: oohama@uec.ac.jp

Abstract—The form of Dueck and Körner’s exponent function for correct decoding probability for discrete memoryless channels at rates above the capacity is similar to the form of Csiszár and Körner’s exponent function for correct decoding probability in lossy source coding for discrete memoryless sources at rates below the rate distortion function. We recently gave a new algorithm for computing Dueck and Körner’s exponent. In this paper, we give an algorithm for computing Csiszár and Körner’s exponent. The proposed algorithm can also be used to compute cutoff rate and the rate distortion function.

Keywords—discrete memoryless source, strong converse, correct decoding probability exponent, iterative algorithm

I. INTRODUCTION

Computation of the channel capacity of a discrete memoryless channel (DMC) under input constraint and computation of the rate distortion function of a discrete memoryless source (DMS) have similar structures [1][2, Chapter 8]. Algorithms for computing channel capacity were given by Blahut [1] and Arimoto [3] and an algorithm for computing rate distortion function was given by Blahut [1].

For channel coding, a strong converse theorem was established by Wolfowitz [4]. Arimoto proved that the probability of correct decoding vanishes exponentially if the transmission rate is above the capacity [5]. He then gave an algorithm for computing his exponent function [6]. Subsequently, Dueck and Körner gave the optimal exponent function of correct decoding probability [7]. They claimed that their exponent function coincides with Arimoto’s exponent. However, the forms of these two exponent functions are quite different. We recently proposed an algorithm for computing Dueck and Körner’s exponent function [8]. The difference between Arimoto’s algorithm and the recently proposed one is as follows: In Arimoto’s algorithm, the probability distribution over the input alphabet and backward transition distribution are updated alternately. On the other hand, in the proposed method, joint probability distribution over input and output alphabets is iteratively updated.

For source coding for DMSs, the rate distribution function, denoted by $R(\Delta|P)$, indicates the minimum admissible rate at distortion level Δ for a source with distribution P . The source coding theorem under ϵ -fidelity criterion states that if the coding rate R is above $R(\Delta|P)$, the probability of an event in which the distortion measure between input sequence and

its reproduced one exceeds Δ tends to zero exponentially [9], [10]. In [6], an algorithm for computing an exponent function of this probability has also been given. On the other hand, the strong converse theorem states that the probability of an event in which the distortion measure exceeds Δ tends to one if $R < R(\Delta|P)$. The optimal exponent for $R < R(\Delta|P)$ was determined by Csiszár and Körner [2]. This exponent function is expressed by a form similar to the form of Dueck and Körner’s exponent function for channel coding. An algorithm for computing the exponent of correct decoding probability for the rates $R < R(\Delta|P)$ has not been provided.

In this paper, we give an iterative algorithm for computing Csiszár and Körner’s exponent function. The algorithm has a structure similar to our recently proposed algorithm for computing Dueck and Körner’s exponent function [8]. We give a proof in which the probability distribution computed by the algorithm converges to the optimal distribution. We also show that the proposed algorithm can be used to compute cutoff rate and the rate distortion function.

Developing a new algorithm for computing the correct decoding probability exponent in lossy source coding has a limited practical importance because the strong converse theorem already states that correct decoding probability goes to zero if the coding rate is below the rate distortion function. The correct decoding exponent expresses how fast such a probability goes to zero. However, analyzing the correct decoding probability exponent and comparing it with the error exponent in source coding as well as the one in channel coding brings a better understanding of the structure of these exponent functions. In addition, the results of this paper may lead to the development of a computation algorithm for other coding schemes.

II. SOURCE CODING AT RATES BELOW THE RATE-DISTORTION FUNCTION

This section gives definitions for quantities that are necessary to describe the correct decoding probability exponent of source coding for discrete memoryless sources (DMSs) at rates below the rate distortion functions.

Let \mathcal{X} be a source alphabet and \mathcal{Y} be a reproduction alphabet. Both \mathcal{X} and \mathcal{Y} are supposed to be finite. A k -length block code for sources with alphabet \mathcal{X} is a pair of mappings $(\varphi^{(k)}, \psi^{(k)})$, where $\varphi^{(k)}$ is an encoding function

that maps every element of \mathcal{X}^k into $\mathcal{M}_k = \{1, 2, \dots, |\mathcal{M}_k|\}$ in a one-to-one manner and $\psi^{(k)}$ is a decoding function that maps every element of \mathcal{M}_k into \mathcal{Y}^k , where \mathcal{M}_k is an index set. The rate of such a code is defined as $\frac{1}{k} \log |\mathcal{M}_k|$. Let $d(x, y) \geq 0$ be a distortion measure for $x \in \mathcal{X}$ and $y \in \mathcal{Y}$. The average distortion between x^k and y^k is defined as $d(x^k, y^k) = \frac{1}{k} \sum_{i=1}^k d(x_i, y_i)$. We assume that for every $x \in \mathcal{X}$, there exists at least one $y \in \mathcal{Y}$ such that $d(x, y) = 0$.

Let P be a probability distribution over source alphabet \mathcal{X} . Correct decoding is defined as an event in which the distortion does not exceed a prescribed distortion level Δ . We denote the probability of correct decoding by $P_c^{(k)}(\varphi^{(k)}, \psi^{(k)}; \Delta|P)$, which is given by

$$P_c^{(k)}(\varphi^{(k)}, \psi^{(k)}; \Delta|P) = \Pr\{d(X^k, \psi^{(k)}(\varphi^{(k)}(X^k))) \leq \Delta\}.$$

The exponent of the maximum of $P_c^{(k)}(\varphi^{(k)}, \psi^{(k)}; \Delta|P)$ over all pairs of encoding and decoding functions having a rate less than R is defined by

$$\begin{aligned} G^{(k)}(R, \Delta|P) \\ := \min_{(\varphi^{(k)}, \psi^{(k)}): \frac{1}{k} \log |\mathcal{M}_k| \leq R} \left(-\frac{1}{k}\right) \log P_c^{(k)}(\varphi^{(k)}, \psi^{(k)}; \Delta|P). \end{aligned}$$

Let

$$G^*(R, \Delta|P) = \lim_{k \rightarrow \infty} G^{(k)}(R, \Delta|P).$$

The optimal exponent $G^*(R, \Delta|P)$ is determined by Csiszár and Körner in [2, p.139]. In order to describe their result, we define

$$\begin{aligned} G_{\text{CK}}(R, \Delta|P) \\ := \min_{q_X \in \mathcal{P}(\mathcal{X})} \{|R(\Delta|q_X) - R|^+ + D(q_X||P)\}, \end{aligned} \quad (1)$$

where $\mathcal{P}(\mathcal{X})$ is a set of probability distributions on \mathcal{X} , $|x|^+ = \max\{0, x\}$,

$$R(\Delta|q_X) = \min_{\substack{q_{Y|X} \in \mathcal{P}(\mathcal{Y}|\mathcal{X}): \\ \mathbb{E}_{q_{XY}}[d(X, Y)] \leq \Delta}} I(q_X, q_{Y|X}),$$

$$D(q_X||P) = \mathbb{E}_{q_X} \left[\log \frac{q_X(X)}{P(X)} \right],$$

where $\mathcal{P}(\mathcal{Y}|\mathcal{X})$ is a set of conditional probability distributions on \mathcal{Y} given \mathcal{X} . Then we have the following theorem:

Theorem 1 (Csiszár and Körner): For any $\Delta \geq 0$ and $0 \leq R \leq R(\Delta|P)$, we have

$$G^*(R, \Delta|P) = G_{\text{CK}}(R, \Delta|P). \quad (2)$$

The purpose of this paper is to give an algorithm for computing $G_{\text{CK}}(R, \Delta|P)$. For this aim, we first introduce the following exponent function and prove it is equivalent to $G_{\text{CK}}(R, \Delta|P)$. We then derive a parametric expression for it. Define

$$\begin{aligned} G(R, \Delta|P) \\ := \min_{\substack{q_{XY} \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}): \\ \mathbb{E}_{q_{XY}}[d(X, Y)] \leq \Delta}} \left\{ |I(q_X, q_{Y|X}) - R|^+ + D(q_X||P) \right\}, \end{aligned} \quad (3)$$

where $\mathcal{P}(\mathcal{X} \times \mathcal{Y})$ is a set of joint distributions on $\mathcal{X} \times \mathcal{Y}$.

We have the following lemma:

Lemma 1: For any $R, \Delta \geq 0$, we have

$$G_{\text{CK}}(R, \Delta|P) = G(R, \Delta|P).$$

Proof: Let q_{XY}^* be a joint distribution that attains $G(R, \Delta|P)$. From its formula, we have

$$R(\Delta|q_X^*) \leq I(q_X^*, q_{Y|X}^*). \quad (4)$$

Thus,

$$\begin{aligned} G(R, \Delta|P) &= |I(q_X^*, q_{Y|X}^*) - R|^+ + D(q_X^*||P) \\ &\stackrel{(a)}{\geq} |R(\Delta|q_X^*) - R|^+ + D(q_X^*||P) \\ &\geq \min_{q_X \in \mathcal{P}(\mathcal{X})} \{|R(\Delta|q_X) - R|^+ + D(q_X||P)\} \\ &= G_{\text{CK}}(R, \Delta|P). \end{aligned}$$

Step (a) follows from (4). On the other hand, let \tilde{q}_X^* be a distribution that attains $G_{\text{CK}}(R, \Delta|P)$ and let $\tilde{q}_{Y|X}^*$ be a conditional distribution that attains $R(\Delta|\tilde{q}_X^*)$. Then, we have

$$\begin{aligned} G_{\text{CK}}(R, \Delta|P) &= |I(\tilde{q}_X^*, \tilde{q}_{Y|X}^*) - R|^+ + D(\tilde{q}_X^*||P) \\ &\geq \min_{\substack{q_{XY}: \\ \mathbb{E}_{q_{XY}}[d(X, Y)] \leq \Delta}} \{|I(q_X, q_{Y|X}) - R|^+ + D(q_X||P)\} \\ &= G(R, \Delta|P). \end{aligned}$$

Thus, we have $G_{\text{CK}}(R, \Delta|P) = G(R, \Delta|P)$, which completes the proof. \blacksquare

The function $G(R, \Delta|P)$ satisfies the following property, which is useful for deriving its parametric expression:

Property 1:

- $G(R, \Delta|P)$ is a monotone decreasing function of $R \geq 0$ for a fixed $\Delta \geq 0$ and is a monotone decreasing function of $\Delta \geq 0$ for a fixed $R \geq 0$.
- $G(R, \Delta|P)$ is a convex function of (R, Δ) .
- $G(R, \Delta|P)$ takes positive value for $0 \leq R < R(\Delta|P)$. For $R \geq R(\Delta|P)$, $G(R, \Delta|P) = 0$.
- For $R' \geq R \geq 0$, we have $G(R, \Delta|P) - G(R', \Delta|P) \leq R' - R$.

See Appendix A for the proof.

In the following, we give definitions of three functions that are related to $G(R, \Delta|P)$ and show their properties. Then we give a lemma, from which a parametric expression of $G(R, \Delta|P)$ is derived.

For $0 \leq \lambda \leq 1, R \geq 0, \Delta \geq 0$ and $\mu \geq 0$, we define

$$\begin{aligned} G^{(\lambda)}(R, \Delta|P) &:= \min_{\substack{q_{XY} \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}): \\ \mathbb{E}_{q_{XY}}[d(X, Y)] \leq \Delta}} \{ \lambda(I(q_X, q_{Y|X}) - R) \\ &\quad + D(q_X||P) \}, \end{aligned} \quad (5)$$

$$\begin{aligned} \Omega^{(\mu, \lambda)}(P) &:= \min_{q_{XY} \in \mathcal{P}(\mathcal{X} \times \mathcal{Y})} \{ \lambda I(q_X, q_{Y|X}) \\ &\quad + D(q_X||P) + \mu \mathbb{E}_{q_{XY}}[d(X, Y)] \}, \end{aligned} \quad (6)$$

$$\begin{aligned} G^{(\mu, \lambda)}(R, \Delta|P) &:= \min_{q_{XY} \in \mathcal{P}(\mathcal{X} \times \mathcal{Y})} \{ \lambda(I(q_X, q_{Y|X}) - R) \\ &\quad + D(q_X||P) - \mu(\Delta - \mathbb{E}_{q_{XY}}[d(X, Y)]) \} \\ &= \Omega^{(\mu, \lambda)}(P) - \lambda R - \mu \Delta. \end{aligned} \quad (7)$$

It is obvious from these definitions that the joint distribution q_{XY} that minimizes $\Omega^{(\mu, \lambda)}$ also minimizes $G^{(\mu, \lambda)}(R, \Delta|P)$ irrespective of the values of R and Δ . $G^{(\lambda)}(R, \Delta|P)$ will be used in Section V for calculating the cut-off rate.

The function $G^{(\lambda)}(R, \Delta|P)$ satisfies the following property:

Property 2:

- a) $G^{(\lambda)}(R, \Delta|P)$ is a monotone decreasing function of $R \geq 0$ for a fixed $\Delta \geq 0$ and is a monotone decreasing function of $\Delta \geq 0$ for a fixed $R \geq 0$.
- b) $G^{(\lambda)}(R, \Delta|P)$ is a convex function of (R, Δ) .

See Appendix B for the proof.

Then, we have the following lemma:

Lemma 2: For any $R \geq 0, \Delta \geq 0$, we have

$$G(R, \Delta|P) = \max_{0 \leq \lambda \leq 1} G^{(\lambda)}(R, \Delta|P). \quad (8)$$

For any $0 \leq \lambda \leq 1, R \geq 0$, and $\Delta \geq 0$, we have

$$G^{(\lambda)}(R, \Delta|P) = \max_{\mu \geq 0} G^{(\mu, \lambda)}(R, \Delta|P). \quad (9)$$

Eqs.(8) and (9) imply that for any $R \geq 0, \Delta \geq 0$, we have

$$G(R, \Delta|P) = \max_{0 \leq \lambda \leq 1} \max_{\mu \geq 0} G^{(\mu, \lambda)}(R, \Delta|P).$$

See Appendix C for the proof. Properties 1 and 2 are needed to prove (8) and (9), respectively.

It follows from Lemma 2 that $G(R, \Delta|P)$ is obtained by maximizing $G^{(\mu, \lambda)}(R, \Delta|P)$ with respect to (μ, λ) . The first step for obtaining $G^{(\mu, \lambda)}(R, \Delta|P)$ is to calculate the joint distribution that minimizes $\Omega^{(\mu, \lambda)}(P)$. In the next section, we give an algorithm to obtain such a joint distribution.

III. DISTRIBUTION UPDATING ALGORITHM

In this section, we propose an iterative algorithm for computing $\Omega^{(\mu, \lambda)}(P)$. Computation of $G(R, \Delta|P)$ from $\Omega^{(\mu, \lambda)}(P)$ is straightforward. We observe that

$$\begin{aligned} \Omega^{(\mu, \lambda)}(P) &= \min_{q_{XY}} \{ \lambda I(q_X, q_{Y|X}) + D(q_X || P) + \mu E_{q_{XY}}[d(X, Y)] \} \\ &= \min_{q_{XY}} E_{q_{XY}} \left[\log \frac{q_X^{1-\lambda}(X) q_{X|Y}^\lambda(X|Y) \exp(\mu d(X, Y))}{P(X)} \right]. \end{aligned}$$

Thus, for computing $\Omega^{(\mu, \lambda)}(P)$, we should find a joint distribution that minimizes the expectation of

$$\omega_q^{(\mu, \lambda)}(x, y) := \log \frac{q_X^{1-\lambda}(x) q_{X|Y}^\lambda(x|y) \exp(\mu d(x, y))}{P(x)}$$

with respect to q_{XY} . Let us define

$$F^{(\mu, \lambda)}(p, q) := E_q \left[\omega_p^{(\mu, \lambda)}(X, Y) \right] + D(q || p),$$

where $p = p_{XY}$ and $q = q_{XY}$ are two probability distributions taking values on $\mathcal{P}(\mathcal{X} \times \mathcal{Y})$.

We have the following two lemmas:

Lemma 3: For a fixed q , $F^{(\mu, \lambda)}(p, q)$ is minimized by $p = q$ and its minimum value is

$$\begin{aligned} F^{(\mu, \lambda)}(q, q) &= E_q[\omega_q^{(\mu, \lambda)}(X, Y)] \\ &= \lambda I(q_X, q_{Y|X}) + D(q_X || P) + \mu E_q[d(X, Y)]. \end{aligned}$$

This implies that

$$\begin{aligned} \min_{p, q} F^{(\mu, \lambda)}(p, q) &= \min_q F^{(\mu, \lambda)}(q, q) \\ &= \min_q E_q[\omega_q^{(\mu, \lambda)}(X, Y)] = \Omega^{(\mu, \lambda)}(P). \end{aligned} \quad (10)$$

Proof: We have

$$\begin{aligned} &F^{(\mu, \lambda)}(p, q) \\ &= E_q \left[\log \frac{p_X^{1-\lambda}(X) p_{X|Y}^\lambda(X|Y)}{P(X) \exp(-\mu d(X, Y))} \right] + E_q \left[\log \frac{q(X, Y)}{p(X, Y)} \right] \\ &= E_q \left[\log \frac{q_X^{1-\lambda}(X) q_{X|Y}^\lambda(X|Y)}{P(X) \exp(-\mu d(X, Y))} \right] \\ &\quad + E_q \left[\log \frac{p_X^{1-\lambda}(X) p_{X|Y}^\lambda(X|Y) q(X, Y)}{q_X^{1-\lambda}(X) q_{X|Y}^\lambda(X|Y) p(X, Y)} \right] \\ &= E_q \left[\omega_q^{(\mu, \lambda)}(X, Y) \right] + E_q \left[\log \frac{q_{Y|X}^{1-\lambda}(Y|X) q_Y^\lambda(Y)}{p_{Y|X}^{1-\lambda}(Y|X) p_Y^\lambda(Y)} \right] \\ &= E_q \left[\omega_q^{(\mu, \lambda)}(X, Y) \right] \\ &\quad + (1 - \lambda) D(q_{Y|X} || p_{Y|X}) + \lambda D(q_Y || p_Y). \end{aligned}$$

Hence, by non-negativity of divergence we have

$$F^{(\mu, \lambda)}(p, q) \geq E_q \left[\omega_q^{(\mu, \lambda)}(X, Y) \right],$$

where equality holds if $p = q$. This completes the proof. ■

Lemma 4: For a fixed p , $F^{(\mu, \lambda)}(p, q)$ is minimized by

$$\begin{aligned} q(x, y) &= \frac{1}{\Lambda_p^{(\mu, \lambda)}} \frac{P(x) \exp(-\mu d(x, y)) p_{XY}(x, y)}{p_X^{1-\lambda}(x) p_{X|Y}^\lambda(x|y)} \\ &:= \hat{q}(p)(x, y), \end{aligned}$$

where $\Lambda_p^{(\mu, \lambda)}$ is a normalization factor defined by

$$\begin{aligned} \Lambda_p^{(\mu, \lambda)} &= E_p \left[\frac{P(X) \exp(-\mu d(X, Y))}{p_X^{1-\lambda}(X) p_{X|Y}^\lambda(X|Y)} \right] \\ &= E_p \left[\exp\{-\omega_p^{(\mu, \lambda)}(X, Y)\} \right] \end{aligned}$$

and its minimum value is

$$\begin{aligned} F^{(\mu, \lambda)}(p, \hat{q}(p)) &= -\log \Lambda_p^{(\mu, \lambda)} \\ &= -\log E_p \left[\exp\{-\omega_p^{(\mu, \lambda)}(X, Y)\} \right]. \end{aligned}$$

This implies that

$$\begin{aligned} \min_{p, q} F^{(\mu, \lambda)}(p, q) &= \min_p F^{(\mu, \lambda)}(p, \hat{q}(p)) \\ &= \min_p \left(-\log E_p \left[\exp\{-\omega_p^{(\mu, \lambda)}(X, Y)\} \right] \right). \end{aligned} \quad (11)$$

Proof: We have

$$\begin{aligned} &F^{(\mu, \lambda)}(p, q) \\ &= E_q \left[\log \frac{p_X^{1-\lambda}(X) p_{X|Y}^\lambda(X|Y) \exp\{\mu d(X, Y)\} q(X, Y)}{P(X) p(X, Y)} \right] \\ &= E_q \left[\log \frac{q(X, Y)}{\frac{1}{\Lambda_p^{(\mu, \lambda)}} \frac{P(X) \exp(-\mu d(X, Y)) p(X, Y)}{p_X^{1-\lambda}(X) p_{X|Y}^\lambda(X|Y)}} \right] - \log \Lambda_p^{(\mu, \lambda)} \\ &\geq -\log \Lambda_p^{(\mu, \lambda)}, \end{aligned}$$

where the last inequality comes from the non-negativity of the divergence. Equality holds if $q(x, y) = \hat{q}(p)(x, y)$. This completes the proof. ■

By Lemmas 3 and 4, we can obtain an iterative algorithm for computing $\Omega^{(\mu, \lambda)}(P)$ as follows:

Distribution updating algorithm

- 1) Choose an initial probability vector $q^{[1]}$ arbitrarily such that all its components are nonzero.
- 2) Then, iterate the following steps for $t = 1, 2, 3, \dots$,

$$q^{[t+1]}(x, y) = \frac{\exp\left\{-\omega_{q^{[t]}}^{(\mu, \lambda)}(x, y)\right\} q^{[t]}(x, y)}{\Lambda_{q^{[t]}}^{(\mu, \lambda)}} \quad (12)$$

$$:= \hat{q}(q^{[t]})(x, y),$$

$$\text{where } \Lambda_{q^{[t]}}^{(\mu, \lambda)} = \mathbb{E}_{q^{[t]}} \left[\exp\left\{-\omega_{q^{[t]}}^{(\mu, \lambda)}(X, Y)\right\} \right].$$

We have the following proposition:

Proposition 1: For $t = 1, 2, 3, \dots$, we have

$$\begin{aligned} F(q^{[1]}, q^{[1]}) &\stackrel{(a)}{\geq} F(q^{[1]}, q^{[2]}) \stackrel{(b)}{\geq} F(q^{[2]}, q^{[2]}) \geq \dots \\ &\geq F(q^{[t]}, q^{[t]}) \\ &\stackrel{(a)}{\geq} F(q^{[t]}, q^{[t+1]}) = -\log \mathbb{E}_{q^{[t]}} \left[\exp\left\{-\omega_{q^{[t]}}^{(\mu, \lambda)}(X, Y)\right\} \right] \\ &\stackrel{(b)}{\geq} F(q^{[t+1]}, q^{[t+1]}) = \mathbb{E}_{q^{[t+1]}} \left[\omega_{q^{[t+1]}}^{(\mu, \lambda)}(X, Y) \right] \\ &\geq \dots \geq \min_q \left\{ -\log \mathbb{E}_q \left[\exp\left\{-\omega_q^{(\mu, \lambda)}(X, Y)\right\} \right] \right\} \\ &\stackrel{(c)}{=} \min_q \mathbb{E}_q \left[\omega_q^{(\mu, \lambda)}(X, Y) \right] = \Omega^{(\mu, \lambda)}(P). \end{aligned}$$

Proof: Step (a) follows from Lemma 4. Step (b) follows from Lemma 3. Step (c) follows from Eq. (10) in Lemma 3 and Eq. (11) in Lemma 4. This completes the proof. ■

IV. CONVERGENCE OF THE ALGORITHM

Proposition 1 shows that $F(q^{[t]}, q^{[t]})$ decreases by updating the probability distribution $q^{[t]}$ according to (12). This section shows that $q^{[t]}$ converges to the optimal distribution. We have the following theorem:

Theorem 2: For any $0 \leq \lambda \leq 1$ and any $\mu \geq 0$ probability vector $q^{[t]}$ defined by (12) converges to the optimal distribution q^* that attains the minimum of

$$\mathbb{E}_q \left[\omega_q^{(\mu, \lambda)}(X, Y) \right] = \lambda I(q_X, q_{Y|X}) + D(q_X \| P) + \mu \mathbb{E}_q[d(X, Y)]$$

in the definition of $\Omega^{(\mu, \lambda)}(P)$.

Proof: By definition, we have $\mathbb{E}_{q^*}[\omega_{q^*}^{(\mu, \lambda)}(X, Y)] = \Omega^{(\mu, \lambda)}(P)$ and $F(q^{[t]}, q^{[t+1]}) = -\log \Lambda_{q^{[t]}}^{(\mu, \lambda)}$. From Eq.(12), we have

$$\Lambda_{q^{[t]}}^{(\mu, \lambda)} = \frac{q^{[t]}(x, y)}{q^{[t+1]}(x, y)} \exp\{-\omega_{q^{[t]}}^{(\mu, \lambda)}(x, y)\}. \quad (13)$$

Hence,

$$\begin{aligned} &-\log \Lambda_{q^{[t]}}^{(\mu, \lambda)} - \Omega^{(\mu, \lambda)}(P) \\ &= -\mathbb{E}_{q^*}[\log \Lambda_{q^{[t]}}^{(\mu, \lambda)}] - \mathbb{E}_{q^*}[\omega_{q^*}^{(\mu, \lambda)}(X, Y)] \\ &\stackrel{(a)}{=} \mathbb{E}_{q^*} \left[\log \frac{q^{[t+1]}(X, Y)}{q^{[t]}(X, Y)} + \omega_{q^{[t]}}^{(\mu, \lambda)}(X, Y) - \omega_{q^*}^{(\mu, \lambda)}(X, Y) \right] \\ &= \mathbb{E}_{q^*} \left[\log \frac{q^{[t+1]}(X, Y)}{q^{[t]}(X, Y)} + \log \left\{ \frac{q_X^{[t]}(X)}{q_X^*(X)} \right\}^{1-\lambda} \right. \\ &\quad \left. + \log \left\{ \frac{q_{Y|X}^{[t]}(X, Y)}{q_{Y|X}^*(X, Y)} \right\}^\lambda \right] \\ &= \mathbb{E}_{q^*} \left[\log \frac{q^{[t+1]}(X, Y)}{q^{[t]}(X, Y)} \right] - (1-\lambda)D(q_X^* \| q_X^{[t]}) \\ &\quad - \lambda D(q_{X|Y}^* \| q_{X|Y}^{[t]} | q_{Y^*}^*) \leq \mathbb{E}_{q^*} \left[\log \frac{q^{[t+1]}(X, Y)}{q^{[t]}(X, Y)} \right], \end{aligned}$$

where equality (a) holds because Eq.(13) holds for all $(x, y) \in \mathcal{X} \times \mathcal{Y}$. Thus,

$$0 \leq -\log \Lambda_{q^{[t]}}^{(\mu, \lambda)} - \Omega^{(\mu, \lambda)}(P) \leq \mathbb{E}_{q^*} \left[\log \frac{q^{[t+1]}(X, Y)}{q^{[t]}(X, Y)} \right] = D(q^* \| q^{[t]}) - D(q^* \| q^{[t+1]}).$$

Therefore, we have $D(q^* \| q^{[t]}) \geq D(q^* \| q^{[t+1]})$, which implies that the KL distance between $q^{[t]}$ and q^* decreases by updating $q^{[t]}$. Put $-\log \Lambda_{q^{[t]}}^{(\mu, \lambda)} - \Omega^{(\mu, \lambda)}(P) = \xi_t$. Then

$$0 \leq \sum_{t=1}^T \xi_t = D(q^* \| q^{[1]}) - D(q^* \| q^{[T+1]}) < D(q^* \| q^{[1]}). \quad (14)$$

$D(q^* \| q^{[1]})$ is finite because all components of $q^{[1]}$ are nonzero. By Proposition 1, $\{\xi_t\}_{t \geq 1}$ is a monotone decreasing sequence. Then from Eq.(14), we have $0 \leq T\xi_T \leq D(q^* \| q^{[1]})$, from which we have

$$0 \leq \xi_T \leq \frac{D(q^* \| q^{[1]})}{T} \rightarrow 0, \quad T \rightarrow \infty.$$

Hence, we have

$$\lim_{t \rightarrow \infty} \left\{ -\log \Lambda_{q^{[t]}}^{(\mu, \lambda)} \right\} = \Omega^{(\mu, \lambda)}(P),$$

completing the proof. ■

As a corollary of Proposition 1 and Theorem 2, we have the following result, which provides a new parametric expression of $G^*(R, \Delta|P) = G(R, \Delta|P)$.

Corollary 1:

$$\begin{aligned} G^{(\mu, \lambda)}(R, \Delta|P) &= -\lambda R - \mu \Delta \\ &\quad + \min_p \left\{ -\log \mathbb{E}_p \left[\frac{P(X) \exp(-\mu d(X, Y))}{p_X^{1-\lambda}(X) p_{Y|X}^\lambda(Y|X)} \right] \right\}, \\ G^*(R, \Delta|P) &= G(R, \Delta|P) \\ &= \max_{0 \leq \lambda \leq 1} \max_{\mu \geq 0} G^{(\mu, \lambda)}(R, \Delta|P) \\ &= \max_{0 \leq \lambda \leq 1} \max_{\mu \geq 0} \min_p \left\{ -\lambda R - \mu \Delta \right. \\ &\quad \left. - \log \mathbb{E}_p \left[\frac{P(X) \exp(-\mu d(X, Y))}{p_X^{1-\lambda}(X) p_{Y|X}^\lambda(Y|X)} \right] \right\}. \end{aligned}$$

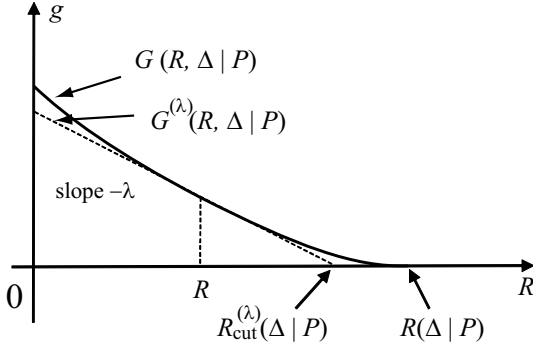


Fig. 1. The bold curve shows the exponent function $G(R, \Delta|P)$ and the dashed line shows the supporting line of slope $-1 \leq -\lambda \leq 0$ to the curve $G(R, \Delta|P)$. $R_{\text{cut}}^{(\lambda)}(\Delta|P)$ is the R -axis intercept of the supporting line, which approaches $R(\Delta|P)$ as $\lambda \rightarrow 0+$.

The proposed algorithm calculates Csiszár and Körner's exponent that expresses the optimal exponent of *correct decoding probability* for $R < R(\Delta|P)$, while Arimoto [6] has presented an iterative algorithm for computing an exponent function of *error probability*¹ derived by Blahut [9]. In Arimoto's algorithm, output distribution $q_Y \in \mathcal{P}(\mathcal{Y})$ and conditional probability distribution $q_{Y|X} \in \mathcal{P}(\mathcal{Y}|\mathcal{X})$ are alternately updated. Unlike Arimoto's algorithm, a joint distribution over the input and output alphabets is updated iteratively in the proposed method. Unfortunately, the proposed algorithm cannot be directly applied to the computations of the exponent function of *error probability* because they involve mini-max structure, *i.e.*, maximization with respect to stochastic matrices and minimization with respect to input distribution.

V. COMPUTATION OF CUTOFF RATE AND THE RATE DISTORTION FUNCTION

The proposed algorithm can be used for computing cutoff rate and the rate distortion function. First, we give the definition of the cutoff rate for lossy source coding. From (5), for a strictly positive λ , we have

$$G^{(\lambda)}(R, \Delta|P) = -\lambda R + \lambda \min_{\substack{q_{XY} \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}): \\ \mathbb{E}_q[d(X, Y)] \leq \Delta}} \{I(q_X, q_{Y|X}) + \frac{1}{\lambda} D(q_X||P)\}. \quad (15)$$

For fixed $\Delta \geq 0$ and $\lambda > 0$, the right hand side of Eq.(15) is viewed as a linear function of R . Moreover, from Eq.(8) in Lemma 2, $G^{(\lambda)}(R, \Delta|P)$ can be viewed as a supporting line to the curve $G(R, \Delta|P)$ with slope $-\lambda$. A rough sketch of the graph $g = G(R, \Delta|P)$ and $g = G^{(\lambda)}(R, \Delta|P)$ is illustrated in Fig. 1. From Property 1, $G(R, \Delta|P)$ takes positive value when $R < R(\Delta|P)$. The cutoff rate is defined as R that satisfies $G^{(\lambda)}(R, \Delta|P) = 0$, *i.e.*,

$$R_{\text{cut}}^{(\lambda)}(\Delta|P) := \min_{\substack{q_{XY} \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}): \\ \mathbb{E}_q[d(X, Y)] \leq \Delta}} \{I(q_X, q_{Y|X}) + \frac{1}{\lambda} D(q_X||P)\}. \quad (16)$$

¹ Blahut [9] gave a lower bound of the exponent function of error probability for $R > R(\Delta|P)$. The optimal exponent of the error probability for $R > R(\Delta|P)$ was determined by Marton [10].

The cutoff rate is calculated by using the proposed method as follows: From Eq.(9) in Lemma 2, we have

$$G^{(\lambda)}(R, \Delta|P) = \max_{\mu \geq 0} G^{(\mu, \lambda)}(R, \Delta|P) = -\lambda R + \lambda \max_{\mu \geq 0} \left\{ \frac{1}{\lambda} \Omega^{(\mu, \lambda)}(P) - \frac{\mu}{\lambda} \Delta \right\}. \quad (17)$$

From Eqs. (15), (16), and (17), we have

$$R_{\text{cut}}^{(\lambda)}(\Delta|P) = \max_{\mu \geq 0} \left\{ \frac{1}{\lambda} \Omega^{(\mu, \lambda)}(P) - \frac{\mu}{\lambda} \Delta \right\}.$$

By Theorem 3, we can calculate $\Omega^{(\mu, \lambda)}(P)$ by the proposed algorithm. Then, the cutoff rate is obtained by calculating the maximum of $\frac{1}{\lambda} \Omega^{(\mu, \lambda)}(P) - \frac{\mu}{\lambda} \Delta$ with respect to $\mu \geq 0$.

We show that $R_{\text{cut}}^{(\lambda)}(\Delta|P)$ approaches $R(\Delta|P)$ as $\lambda \rightarrow 0+$. Let $\alpha = \min\{\log|\mathcal{X}|, \log|\mathcal{Y}|\}$ and $d_{\text{max}} = \max_{(x, y) \in (\mathcal{X}, \mathcal{Y})} d(x, y)$. We have the following proposition.

Proposition 2: $R_{\text{cut}}^{(\lambda)}(\Delta|P)$ is a monotone decreasing function of λ . Moreover, if $\lambda \leq \frac{1}{8\alpha}$, we have

$$0 \leq R(\Delta|P) - R_{\text{cut}}^{(\lambda)}(\Delta|P) \leq c_1 \sqrt{\lambda} (\log \lambda^{-1} + c_2), \quad (18)$$

where $c_1 = \frac{3}{2} \sqrt{2\alpha}$, $c_2 = \frac{4}{3} \log(|\mathcal{X}||\mathcal{Y}|) - \log(2\alpha) + \frac{2}{3} d_{\text{max}} |R'(\Delta|P)|$, and $R'(\Delta|P) = \frac{d}{d\Delta} R(\Delta|P)$. This inequality implies that

$$\lim_{\lambda \rightarrow 0+} R_{\text{cut}}^{(\lambda)}(\Delta|P) = R(\Delta|P).$$

See Appendix D for the proof.

This proposition implies that by choosing a sufficiently small $\lambda > 0$, we can use $R_{\text{cut}}^{(\lambda)}(\Delta|P)$ as a good approximation of $R(\Delta|P)$ for which accuracy is guaranteed by (18).

VI. CONCLUSION

We have proposed an iterative algorithm for computing Csiszár and Körner's exponent [2] that expresses the optimal exponent of correct decoding probability in lossy source coding when a rate R is below the rate distortion function $R(\Delta|P)$. The proposed algorithm has a structure similar to the one proposed by the authors [8] that computes Dueck and Körner's exponent in channel coding when the rate is above the capacity. We have proven the joint distribution calculated by the proposed algorithm converges to the optimal distribution that achieves Csiszár and Körner's exponent. We have also shown that our proposed algorithm can be used to calculate cutoff rate and the rate distortion function.

ACKNOWLEDGMENT

This work was supported in part by JSPS KAKENHI Grant Numbers 23360172, 25820162 and K16000333.

REFERENCES

- [1] R. E. Blahut, "Computation of channel capacity and rate-distortion functions," *IEEE Trans. Inform. Theory*, vol. 18, pp. 460–473, 1972.
- [2] I. Csiszár and J. Körner, *Information theory, coding theorems for discrete memoryless systems*. Cambridge University Press, 2nd edition, 2011.
- [3] S. Arimoto, "An algorithm for computing the capacity of arbitrary discrete memoryless channels," *IEEE Trans. Inform. Theory*, vol. 18, pp. 14–20, 1972.
- [4] J. Wolfowitz, "Products of indecomposable, aperiodic, stochastic matrices," *Proc. Amer. Math. Soc.*, vol. 14, no. 5, pp. 733–737, 1963.
- [5] S. Arimoto, "On the converse to the coding theorem for discrete memoryless channels," *IEEE Trans. Inform. Theory*, vol. 19, pp. 357–359, 1973.
- [6] S. Arimoto, "Computation of random coding exponent functions," *IEEE Trans. Inform. Theory*, vol. 22, pp. 665–671, 1976.
- [7] G. Dueck and J. Körner, "Reliability function of a discrete memoryless channel at rates above capacity," *IEEE Trans. Inform. Theory*, vol. 25, pp. 82–85, 1979.
- [8] Y. Oohama and Y. Jitsumatsu, "A new iterative algorithm for computing the optimal exponent of correct decoding for discrete memoryless channels," *IEEE Int. Symp. Inform. Theory (ISIT2015)*, 2015.
- [9] R. E. Blahut, "Hypothesis testing and information theory," *IEEE Trans. Inform. Theory*, vol. 20, pp. 405–417, 1974.
- [10] D. R. Marton, "Error exponent for source coding with a fidelity criterion," *IEEE Trans. Inform. Theory*, vol. 20, pp. 197–199, 1974.

APPENDIX

A. Proof of Property 1

In this appendix, we prove Property 1. By definition, Part a) is obvious. For the proof of Part b), let $q^{(0)}$ and $q^{(1)}$ be joint distribution functions that attain $G(R_0, \Delta_0|P)$ and $G(R_1, \Delta_1|P)$, respectively. Denote

$$\begin{aligned} \Theta(R, q|P) &:= |I(q_X, q_{Y|X}) - R|^+ + D(q_X||P) \\ &= \max\{D(q_X||P), I(q_X, q_{Y|X}) - R + D(q_X||P)\}. \end{aligned} \quad (19)$$

By definition, we have

$$G(R_i, \Delta_i|P) = \Theta(R_i, q^{(i)}|P) \text{ for } i = 0, 1. \quad (20)$$

For $\alpha_1 = \alpha \in [0, 1]$ and $\alpha_0 = 1 - \alpha$, we set $R_\alpha = \alpha_0 R_0 + \alpha_1 R_1$, $\Delta_\alpha = \alpha_0 \Delta_0 + \alpha_1 \Delta_1$, and $q^{(\alpha)} = \alpha_0 q^{(0)} + \alpha_1 q^{(1)}$. By linearity of $E_q[d(X, Y)]$ with respect to q , we have that

$$E_{q^{(\alpha)}}[d(X, Y)] = \sum_{i=0,1} \alpha_i E_{q^{(i)}}[d(X, Y)] \leq \Delta_\alpha. \quad (21)$$

Because

$$I(q_X, q_{Y|X}) + D(q_X||P) = \sum_{x,y} q_{XY}(x, y) \log \frac{q_{X|Y}(x|y)}{P(x)}$$

is convex with respect to q_{XY} and $D(q_X||P)$ is convex with respect to q_X , we have

$$\begin{aligned} &I(q_X^{(\alpha)}, q_{Y|X}^{(\alpha)}) + D(q_X^{(\alpha)}||P) \\ &\leq \sum_{i=0,1} \alpha_i \left\{ I(q_X^{(i)}, q_{Y|X}^{(i)}) + D(q_X^{(i)}||P) \right\}, \end{aligned} \quad (22)$$

$$D(q_X^{(\alpha)}||P) \leq \sum_{i=0,1} \alpha_i D(q_X^{(i)}||P). \quad (23)$$

Therefore, we have the following two chains of inequalities:

$$\begin{aligned} &I(q_X^{(\alpha)}, q_{Y|X}^{(\alpha)}) + D(q_X^{(\alpha)}||P) - R_\alpha \\ &\stackrel{(a)}{\leq} \sum_{i=0,1} \alpha_i \left\{ I(q_X^{(i)}, q_{Y|X}^{(i)}) + D(q_X^{(i)}||P) - R_i \right\} \\ &\stackrel{(b)}{\leq} \sum_{i=0,1} \alpha_i \Theta(R_i, q^{(i)}|P), \end{aligned} \quad (24)$$

$$\begin{aligned} &D(q_X^{(\alpha)}||P) \stackrel{(c)}{\leq} \sum_{i=0,1} \alpha_i D(q_X^{(i)}||P) \\ &\stackrel{(d)}{\leq} \sum_{i=0,1} \alpha_i \Theta(R_i, q^{(i)}|P). \end{aligned} \quad (25)$$

Steps (a) and (c) follow from (22) and (23) and Steps (b) and (d) follow from the definition of $\Theta(R_i, q^{(i)}|P)$ for $i = 0, 1$. Then, from (19) we have

$$\Theta(R_\alpha, q^{(\alpha)}|P) \leq \sum_{i=0,1} \alpha_i \Theta(R_i, q^{(i)}|P). \quad (26)$$

Therefore,

$$\begin{aligned} G(R_\alpha, \Delta_\alpha|P) &= \min_{\substack{q \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}): \\ E_q[d(X, Y)] \leq \Delta_\alpha}} \Theta(R_\alpha, q|P) \\ &\stackrel{(a)}{\leq} \Theta(R_\alpha, q^{(\alpha)}|P) \stackrel{(b)}{\leq} \sum_{i=0,1} \alpha_i \Theta(R_i, q^{(i)}|P) \\ &\stackrel{(c)}{=} \sum_{i=0,1} \alpha_i G(R_i, \Delta_i|P). \end{aligned}$$

Step (a) follows from (21), Step (b) follows from (26), and Step (c) follows from (20).

For the proof of Part c), the choice of $q_X = P$ gives $G(R, \Delta|P) = 0$, if $R \geq R(\Delta|P)$. If $R < R(\Delta|P)$, the choice of $q_X = P$ makes the first term of the objective function strictly positive, while any choice of $q \neq P$, $D(q||P)$ is strictly positive. This completes the proof of Part c).

For the proof of Part d), let q^* be a joint distribution that attains $G(R', \Delta|P)$. Then,

$$\begin{aligned} G(R, \Delta|P) &\leq |I(q_X^*, q_{Y|X}^*) - R|^+ + D(q_X^*||P) \\ &\stackrel{(a)}{\leq} (R' - R) + |I(q_X^*, q_{Y|X}^*) - R'|^+ \\ &\quad + D(q_X^*||P) \\ &= (R' - R) + G(R', \Delta|P). \end{aligned}$$

Step (a) follows from $|x|^+ \leq |x - c|^+ + c$ for $c \geq 0$. This completes the proof. \blacksquare

B. Proof of Property 2

In this appendix we prove Property 2. By definition, Part a) is obvious. For the proof of Part b), let $q^{(0)}$ and $q^{(1)}$ be joint distribution functions that attain $G^{(\lambda)}(R_0, \Delta_0|P)$ and $G^{(\lambda)}(R_1, \Delta_1|P)$, respectively. Denote

$$\Theta^{(\lambda)}(R, q|P) := \lambda[I(q_X, q_{Y|X}) - R] + D(q_X||P). \quad (27)$$

By definition, we have

$$G^{(\lambda)}(R_i, \Delta_i|P) = \Theta^{(\lambda)}(R_i, q^{(i)}|P) \text{ for } i = 0, 1. \quad (28)$$

For $\alpha_1 = \alpha \in [0, 1]$ and $\alpha_0 = 1 - \alpha$, we set $R_\alpha = \alpha_0 R_0 + \alpha_1 R_1$, $\Delta_\alpha = \alpha_0 \Delta_0 + \alpha_1 \Delta_1$, and $q^{(\alpha)} = \alpha_0 q^{(0)} + \alpha_1 q^{(1)}$. By linearity of $\mathbb{E}_q[d(X, Y)]$ with respect to q , we have that

$$\mathbb{E}_{q^{(\alpha)}}[d(X, Y)] = \sum_{i=0,1} \alpha_i \mathbb{E}_{q^{(i)}}[d(X, Y)] \leq \Delta_\alpha. \quad (29)$$

Since

$$\begin{aligned} & \lambda I(q_X, q_{Y|X}) + D(q_X \| P) \\ &= \lambda [I(q_X, q_{Y|X}) + D(q_X \| P)] + (1 - \lambda) D(q_X \| P) \\ &= \lambda \sum_{x,y} q_{XY}(x, y) \log \frac{q_{X|Y}(x|y)}{P(x)} + (1 - \lambda) D(q_X \| P) \end{aligned}$$

is convex with respect to q_{XY} , we have

$$\begin{aligned} & \Theta^{(\lambda)}(R_\alpha, q^{(\alpha)} | P) \\ &= \lambda [I(q_X^{(\alpha)}, q_{Y|X}^{(\alpha)}) - R_\alpha] + D(q_X^{(\alpha)} \| P) \\ &\leq \sum_{i=0,1} \alpha_i \left\{ \lambda [I(q_X^{(i)}, q_{Y|X}^{(i)}) - R_i] + D(q_X^{(i)} \| P) \right\} \\ &\stackrel{(a)}{=} \sum_{i=0,1} \alpha_i \Theta^{(\lambda)}(R_i, q^{(i)} | P). \end{aligned} \quad (30)$$

Step (a) follows from the definition of $\Theta^{(\lambda)}(R_i, q^{(i)} | P)$ for $i = 0, 1$. Therefore,

$$\begin{aligned} G^{(\lambda)}(R_\alpha, \Delta_\alpha | P) &= \min_{\substack{q \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}): \\ \mathbb{E}_q[d(X, Y)] \leq \Delta_\alpha}} \Theta^{(\lambda)}(R_\alpha, q | P) \\ &\stackrel{(a)}{\leq} \Theta^{(\lambda)}(R_\alpha, q^{(\alpha)} | P) \stackrel{(b)}{\leq} \sum_{i=0,1} \alpha_i \Theta^{(\lambda)}(R_i, q^{(i)} | P) \\ &\stackrel{(c)}{=} \sum_{i=0,1} \alpha_i G^{(\lambda)}(R_i, \Delta_i | P). \end{aligned}$$

Step (a) follows from (29), Step (b) follows from (30), and Step (c) follows from (28). This completes the proof. ■

C. Proof of Lemma 2

In this appendix we prove Lemma 2. First, we prove Eq.(8). For any $\lambda \in [0, 1]$, we have $|x|^+ \geq \lambda x$. Let \hat{q} be a joint distribution that attains $G(R, \Delta | P)$. Then, we have

$$\begin{aligned} G(R, \Delta | P) &= |I(\hat{q}_X, \hat{q}_{Y|X}) - R|^+ + D(\hat{q}_X \| P) \\ &\geq \lambda [I(\hat{q}_X, \hat{q}_{Y|X}) - R] + D(\hat{q}_X \| P) \\ &\geq \min_{\substack{q \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}): \\ \mathbb{E}_q[d(X, Y)] \leq \Delta}} \{ \lambda [I(q_X, q_{Y|X}) - R] + D(q_X \| P) \} \\ &= G^{(\lambda)}(R, \Delta | P). \end{aligned}$$

Thus,

$$G(R, \Delta | P) \geq \max_{0 \leq \lambda \leq 1} G^{(\lambda)}(R, \Delta | P).$$

Hence, it is sufficient to show that there exists a $\lambda \in [0, 1]$ such that $G(R, \Delta | P) \leq G^{(\lambda)}(R, \Delta | P)$. From Property 1, there exists a $\lambda \in [0, 1]$ such that for any $R' \geq 0$ we have

$$G(R', \Delta | P) \geq G(R, \Delta | P) - \lambda(R' - R). \quad (31)$$

Fix the above λ . Let q^* be a joint distribution that attains $G^{(\lambda)}(R, \Delta | P)$. Set $R' = I(q_X^*, q_{Y|X}^*)$. Then we have

$$\begin{aligned} & G(R, \Delta | P) \\ &\stackrel{(a)}{\leq} G(R', \Delta | P) + \lambda(R' - R) \\ &= \min_{\substack{q \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}): \\ \mathbb{E}_q[d(X, Y)] \leq \Delta}} \{ |I(q_X, q_{Y|X}) - R'|^+ + D(q_X \| P) \} \\ &\quad + \lambda(R' - R) \\ &\leq |I(q_X^*, q_{Y|X}^*) - R'|^+ + D(q_X^* \| P) + \lambda(R' - R) \\ &\stackrel{(b)}{=} D(q_X^* \| P) + \lambda [I(q_X^*, q_{Y|X}^*) - R] \\ &= G^{(\lambda)}(R, \Delta | P). \end{aligned} \quad (32)$$

Step (a) follows from (31) and Step (b) comes from the choice of $R' = I(q_X^*, q_{Y|X}^*)$. Therefore, there exists a $0 \leq \lambda \leq 1$ such that $G^{(\mu)}(R, \Delta | P) = G^{(\mu, \lambda)}(R, \Delta | P)$.

Next, we prove (9). From its formula, it is obvious that

$$G^{(\lambda)}(R, \Delta | P) \geq \max_{\mu \geq 0} G^{(\mu, \lambda)}(R, \Delta | P).$$

Hence, it is sufficient to show that for any $R \geq 0$ and $\Delta \geq 0$, there exists $\mu \geq 0$ such that

$$G^{(\lambda)}(R, \Delta | P) \leq G^{(\mu, \lambda)}(R, \Delta | P). \quad (33)$$

From Property 2 part a) and b), $G^{(\lambda)}(R, \Delta | P)$ is a monotone decreasing and convex function of $\Delta \geq 0$ for a fixed R . Thus, there exists $\mu \geq 0$ such that for any $\Delta' \geq 0$, the following inequality holds:

$$G^{(\lambda)}(R, \Delta' | P) \geq G^{(\lambda)}(R, \Delta | P) - \mu(\Delta' - \Delta). \quad (34)$$

Fix the above μ . Let q^* be a joint distribution that attains $G^{(\mu, \lambda)}(R, \Delta | P)$. Set $\Delta' = \mathbb{E}_{q^*}[d(X, Y)]$. Then, we have

$$\begin{aligned} G^{(\lambda)}(R, \Delta | P) &\stackrel{(a)}{\leq} G^{(\lambda)}(R, \Delta' | P) - \mu(\Delta - \Delta') \\ &= \min_{\substack{q: \\ \mathbb{E}_q[d(X, Y)] \leq \Delta'}} \{ \lambda [I(q_X, q_{Y|X}) - R] + D(q_X \| P) \} \\ &\quad - \mu(\Delta - \Delta') \\ &\stackrel{(b)}{\leq} \lambda [I(q_X^*, q_{Y|X}^*) - R] + D(q_X^* \| P) - \mu(\Delta - \mathbb{E}_{q^*}[d(X, Y)]) \\ &= G^{(\mu, \lambda)}(R, \Delta | P). \end{aligned}$$

Step (a) follows from (34) and Step (b) follows from the definition of $G^{(\lambda)}(R, \Delta' | P)$ and the choice of $\Delta' = \mathbb{E}_{q^*}[d(X, Y)]$. Thus, for any $\Delta \geq 0$, we have (33) for some $\mu \geq 0$. This completes the proof. ■

D. Proof of Proposition 2

In this appendix, we prove Proposition 2. We begin with the following lemma:

Lemma 5: If two probability distributions p and q on \mathcal{X} satisfy $D(p||q) \leq \nu$ for a constant $\nu \leq \frac{1}{8}$, we have

$$|H(p) - H(q)| \leq \sqrt{2\nu} \log \frac{|\mathcal{X}|}{\sqrt{2\nu}}.$$

Proof: From Pinsker's inequality, we have

$$D(p||q) \geq \frac{1}{2} \|p - q\|_1,$$

where $\|p - q\|_1 = \sum_{x \in \mathcal{X}} |p(x) - q(x)|$. It follows from Lemma 2.7 in [2] that if $\|p - q\|_1 = \Theta \leq \frac{1}{2}$, then we have

$$|H(p) - H(q)| \leq \Theta \log \frac{|\mathcal{X}|}{\Theta}.$$

The lemma is proved by combining these two inequalities together with monotone increasing property of $\Theta \log \frac{|\mathcal{X}|}{\Theta}$ for $0 \leq \Theta \leq |\mathcal{X}|/e$. ■

Proof of Proposition 2: First we show the monotonicity of $R_{\text{cut}}^{(\lambda)}(R, \Delta|P)$ with respect to λ . Let $0 < \lambda \leq \lambda' \leq 1$ and q^* be a joint distribution that attains $R_{\text{cut}}^{(\lambda)}(R, \Delta|P)$. Then, we have

$$\begin{aligned} R_{\text{cut}}^{(\lambda')}(R, \Delta|P) &\stackrel{(a)}{\leq} I(q_X^*, q_{Y|X}^*) + \frac{1}{\lambda'} D(q_X^* \| P) \\ &\stackrel{(b)}{\leq} I(q_X^*, q_{Y|X}^*) + \frac{1}{\lambda} D(q_X^* \| P) \\ &= R_{\text{cut}}^{(\lambda)}(R, \Delta|P). \end{aligned}$$

Step (a) follows from the definition and step (b) follows from $\lambda \leq \lambda'$.

Next, we prove (18). Let V^* be a distribution on \mathcal{Y} given \mathcal{X} that attains $R(\Delta|P)$. Then, the choice of $(q_X, q_{Y|X}) = (P, V^*)$ gives

$$R_{\text{cut}}^{(\lambda)}(\Delta|P) \leq I(P, V^*) = R(\Delta|P). \quad (35)$$

This gives the first inequality in (18).

For the proof of second inequality in (18), we first give an lower bound of $R_{\text{cut}}^{(\lambda)}(\Delta|P)$ and then give an upper bound of $R(\Delta|P)$. Let q^* be a joint distribution that attains $R_{\text{cut}}^{(\lambda)}(R, \Delta|P)$. Then, we have

$$R_{\text{cut}}^{(\lambda)}(\Delta|P) = I(q_X^*, q_{Y|X}^*) + \frac{1}{\lambda} D(q_X^* \| P).$$

By the non-negativity of divergence, we have

$$R_{\text{cut}}^{(\lambda)}(\Delta|P) \geq I(q_X^*, q_{Y|X}^*). \quad (36)$$

By the non-negativity of mutual information, we also have

$$\begin{aligned} \frac{1}{\lambda} D(q_X^* \| P) &\leq R_{\text{cut}}^{(\lambda)}(\Delta|P) \\ &\stackrel{(a)}{\leq} I(P, V^*) \leq \min\{\log |\mathcal{X}|, \log |\mathcal{Y}|\}. \end{aligned}$$

Step (a) follows from (35). Let $\alpha = \min\{\log |\mathcal{X}|, \log |\mathcal{Y}|\}$. Then,

$$D(q_X^* \| P) \leq \alpha \lambda. \quad (37)$$

Thus, $D(q_X^* \| P) \rightarrow 0$ as $\lambda \rightarrow 0+$, which shows q_X^* converges to P . We have

$$\begin{aligned} &I(P, q_{Y|X}^*) - I(q_X^*, q_{Y|X}^*) \\ &= H(P) + H(Pq_{Y|X}^*) - H((P, q_{Y|X}^*)) \\ &\quad - \{H(q_X^*) + H(q_Y^*) - H(q_{XY}^*)\} \\ &\leq |H(P) - H(q_X^*)| + |H(Pq_{Y|X}^*) - H(q_Y^*)| \\ &\quad + |H((P, q_{Y|X}^*)) - H(q_{XY}^*)|. \end{aligned} \quad (38)$$

From Lemma 5 and (37), the first term of (38) is upper bounded by $\sqrt{2\alpha\lambda} \log \frac{|\mathcal{X}|}{\sqrt{2\alpha\lambda}}$ if $\alpha\lambda \leq \frac{1}{8}$. By the chain rule of the divergence, we have $D(q_{XY}^* \| (P, q_{Y|X}^*)) = D(q_X^* \| P) + D(q_{Y|X}^* \| q_X^*) = D(q_X^* \| P)$ and $D(Pq_{Y|X}^* \| q_Y^*) \leq$

$D(q_{XY}^* \| (P, q_{Y|X}^*))$. Thus, the second and the third terms of (38) are upper bounded by $\sqrt{2\alpha\lambda} \log \frac{|\mathcal{Y}|}{\sqrt{2\alpha\lambda}}$ and $\sqrt{2\alpha\lambda} \log \frac{|\mathcal{X}||\mathcal{Y}|}{\sqrt{2\alpha\lambda}}$, respectively. Therefore, by (36) and (38), if $\lambda \leq \frac{1}{8\alpha}$, we have

$$\begin{aligned} R_{\text{cut}}^{(\lambda)}(\Delta|P) &\geq I(q_X^*, q_{Y|X}^*) \\ &\geq I(P, q_{Y|X}^*) - \sqrt{2\alpha\lambda} \log \frac{|\mathcal{X}|}{\sqrt{2\alpha\lambda}} \\ &\quad - \sqrt{2\alpha\lambda} \log \frac{|\mathcal{Y}|}{\sqrt{2\alpha\lambda}} - \sqrt{2\alpha\lambda} \log \frac{|\mathcal{X}||\mathcal{Y}|}{\sqrt{2\alpha\lambda}} \\ &= I(P, q_{Y|X}^*) - \sqrt{2\alpha\lambda} \log \frac{|\mathcal{X}|^2 |\mathcal{Y}|^2}{(2\alpha\lambda)^{\frac{3}{2}}}. \end{aligned} \quad (39)$$

Next, we give an upper bound of $R(\Delta|P)$. By the convexity of $R(\Delta|P)$, for any ν , we have

$$R(\Delta + \nu) \geq R(\Delta|P) + \nu R'(\Delta|P), \quad (40)$$

where

$$R'(\Delta|P) = \frac{d}{d\Delta} R(\Delta|P).$$

Note that $R(\Delta|P)$ is monotone decreasing function of Δ and thus $R'(\Delta|P) \leq 0$. We have

$$\begin{aligned} &E_{(P, q_{Y|X}^*)}[d(X, Y)] - E_{q^*}[d(X, Y)] \\ &= \sum_{x \in \mathcal{X}} [P(x) - q_X^*(x)] \sum_{y \in \mathcal{Y}} q_{Y|X}^*(y|x) d(x, y) \\ &\leq \sum_{x \in \mathcal{X}} |P(x) - q_X^*(x)| d_{\max} = \|P - q_X^*\|_1 d_{\max}. \end{aligned}$$

Therefore, the following inequality holds:

$$\begin{aligned} E_{(P, q_{Y|X}^*)}[d(X, Y)] &\leq E_{q^*}[d(X, Y)] + \|P - q_X^*\|_1 d_{\max} \\ &\stackrel{(a)}{\leq} E_{q^*}[d(X, Y)] + \sqrt{2\alpha\lambda} d_{\max} \\ &\stackrel{(b)}{\leq} \Delta + \sqrt{2\alpha\lambda} d_{\max}. \end{aligned} \quad (41)$$

Step (a) follows from (37) and Pinsker's inequality. Step (b) follows from the definition of q^* . Then, we have

$$\begin{aligned} R(\Delta + \sqrt{2\alpha\lambda} d_{\max}|P) &= \min_{\substack{W \in \mathcal{P}(\mathcal{Y}|\mathcal{X}): \\ E_{(P, W)}[d(X, Y)] \leq \Delta + \sqrt{2\alpha\lambda} d_{\max}}} I(P, W) \\ &\stackrel{(a)}{\leq} I(P, q_{Y|X}^*). \end{aligned} \quad (42)$$

Step (a) follows from (41). Then, we have the following inequality:

$$\begin{aligned} R(\Delta|P) &\stackrel{(a)}{\leq} R(\Delta + \sqrt{2\alpha\lambda} d_{\max}|P) - \sqrt{2\alpha\lambda} d_{\max} R'(\Delta|P) \\ &\stackrel{(b)}{\leq} I(P, q_{Y|X}^*) + \sqrt{2\alpha\lambda} d_{\max} |R'(\Delta|P)|. \end{aligned} \quad (43)$$

Step (a) follows from (40) with $\nu = \sqrt{2\alpha\lambda} d_{\max}$. Step (b) follows from (42).

Then, we have the following:

$$\begin{aligned}
& R(\Delta|P) - R_{\text{cut}}^{(\lambda)}(\Delta|P) \\
& \stackrel{(a)}{\leq} I(P, q_{Y|X}^*) + \sqrt{2\alpha\lambda} d_{\max} |R'(\Delta|P)| \\
& \quad - \{I(P, q_{Y|X}^*) - \sqrt{2\alpha\lambda} \log \frac{|\mathcal{X}|^2 |\mathcal{Y}|^2}{(2\alpha\lambda)^{\frac{3}{2}}}\} \\
& = \sqrt{2\alpha\lambda} \left\{ \log \frac{|\mathcal{X}|^2 |\mathcal{Y}|^2}{(2\alpha\lambda)^{\frac{3}{2}}} + d_{\max} |R'(\Delta|P)| \right\} \\
& = \sqrt{2\alpha\lambda} \left\{ -\frac{3}{2} \log \lambda + \log \frac{|\mathcal{X}|^2 |\mathcal{Y}|^2}{(2\alpha)^{\frac{3}{2}}} + d_{\max} |R'(\Delta|P)| \right\} \\
& = c_1 \sqrt{\lambda} (\log \lambda^{-1} + c_2),
\end{aligned}$$

where $c_1 = \frac{3}{2} \sqrt{2\alpha}$, $c_2 = \frac{4}{3} \log(|\mathcal{X}||\mathcal{Y}|) - \log(2\alpha) + \frac{2}{3} d_{\max} |R'(\Delta|P)|$. Step (a) follows from (39) and (43). This completes the proof. \blacksquare